

Variational Inference For Probabilistic Latent Tensor Factorization with KL Divergence

Beyza Ermiş¹, Y. Kenan Yılmaz¹, A. Taylan Cemgil¹ and Evrim Acar²

¹*Department of Computer Science, Boğaziçi University, Istanbul, Turkey*

²*Faculty of Life Sciences, University of Copenhagen, Frederiksberg C, Denmark*

Abstract—Probabilistic Latent Tensor Factorization (PLTF) is a recently proposed probabilistic framework for modelling multi-way data. Not only the common tensor factorization models but also any arbitrary tensor factorization structure can be realized by the PLTF framework. This paper presents full Bayesian inference via variational Bayes that facilitates more powerful modelling and allows more sophisticated inference on the PLTF framework. We illustrate our approach on model order selection and link prediction.

Keywords—Probabilistic Latent Tensor Factorization(PLTF); Variational Bayes(VB); Link Prediction; missing data

I. INTRODUCTION

Factorization based data modelling has become popular together with the advances in the computational power. Non-negative Matrix Factorization (NMF) model, proposed by Lee and Seung [1] (and also earlier by Paatero and Tapper [2]), is one of the most popular factorization models where the aim is to estimate the matrices Z_1 and Z_2 as the matrix X is observed:

$$X(i, j) \approx \hat{X}(i, j) = \sum_k Z_1(i, k) Z_2(k, j). \quad (1)$$

Here X , Z_1 and Z_2 are all non-negative matrices. This modelling paradigm has found place in many fields including recommender systems [3], image processing [4] and bioinformatics [5].

Although the NMF model has its own advantages, certain applications require more structured modelling and incorporation of prior knowledge where NMF can be inadequate. Accordingly, several complex factorization models have been proposed in the literature [5]. The probabilistic Latent Tensor Factorization framework (PLTF) [6] enables one to incorporate domain specific information to any arbitrary factorization model and provides the update rules for multiplicative gradient descent and expectation-maximization algorithms.

The PLTF framework is defined as a natural extension of the matrix factorization model of (1):

$$X(v_0) \approx \hat{X}(v_0) = \sum_{\bar{v}_0} \prod_{\alpha} Z_{\alpha}(v_{\alpha}), \quad (2)$$

where $\alpha = 1, \dots, K$ denotes the factor index. In this framework, the goal is to compute an approximate factorization

of a given higher-order tensor, i.e., a multiway array, X in terms of a product of individual factors Z_{α} , some of which are possibly fixed. Here, we define V as the set of all indices in a model, V_0 as the set of visible indices, V_{α} as the set of indices in Z_{α} , and $\bar{V}_{\alpha} = V - V_{\alpha}$ as the set of all indices not in Z_{α} . We use small letters as v_{α} to refer to a particular setting of indices in V_{α} . Since the product $\prod_{\alpha} Z_{\alpha}(v_{\alpha})$ is collapsed over a set of indices, the factorization is latent.

In this study, we use non-negative variants of the two most widely-used low-rank tensor factorization models; the Tucker model [7] and the more restricted CANDECOMP/PARAFAC (CP) model [8], [9], [10]. In order to illustrate the approach, we can define these models in the PLTF notation. Given a three-way tensor X the CP model is defined as follows:

$$X(i, j, k) \approx \hat{X}(i, j, k) = \sum_r Z_1(i, r) Z_2(j, r) Z_3(k, r) \quad (3)$$

where the index sets $V = \{i, j, k, r\}$, $V_0 = \{i, j, k\}$, $V_1 = \{i, r\}$, $V_2 = \{j, r\}$ and $V_3 = \{k, r\}$. An alternative Tucker model of X is defined in the PLTF notation as follows:

$$\hat{X}(i, j, k) = \sum_{p, q, r} Z_1(i, p) Z_2(j, q) Z_3(k, r) Z_4(p, q, r) \quad (4)$$

where the index sets $V = \{i, j, k, p, q, r\}$, $V_0 = \{i, j, k\}$, $V_1 = \{i, p\}$, $V_2 = \{j, q\}$, $V_3 = \{k, r\}$ and $V_4 = \{p, q, r\}$.

The main contributions of this paper can be summarized as follows:

- Variational Bayes procedure for making inference on the PLTF framework is presented.
- Exact characterization of the approximating distribution and full conditionals are observed as a product of multinomial distributions, leading to a richer approximation distribution than a naive mean field.
- Computation of a variational lower bound for estimation of marginal likelihood of a tensor factorization model is described.
- A model selection framework for arbitrary non-negative tensor factorization model for KL cost with the variational bound is constructed.

- The proposed approach is illustrated on link prediction problem: the problem of predicting the existence of connections between entities of interest.

A. Probability Model

The usual approach to estimate the factors Z_α is trying to find the optimal $Z_{1:K}^* = \operatorname{argmin}_{Z_{1:K}} d(X||\hat{X})$, where $d(\cdot)$ is a divergence typically taken as Euclidean, Kullback-Leibler or Itakura-Saito divergences. Since the analytical solution for this problem is intractable, one should refer to iterative or approximate inference methods.

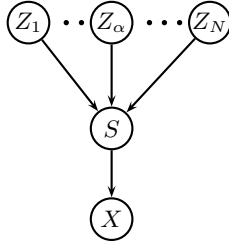


Figure 1. The generative model of the PLTF framework as a Bayesian network. The directed acyclic graph describes the dependency structure of the variables: the full joint distribution can be written as $p(X, S, Z_{1:K}) = p(X|S)p(S|Z_{1:K})\prod_{\alpha} p(Z_{\alpha})$.

In this study, we use the Kullback-Leibler (KL) divergence as the cost function which is equivalent to selecting the Poisson observation model [4], [6], while our approach can be extended to other costs where a composite structure is present. The overall probabilistic model is defined as follows:

$$Z_{\alpha}(v_{\alpha}) \sim \mathcal{G}(Z_{\alpha}(v_{\alpha}); A_{\alpha}(v_{\alpha}), B_{\alpha}(v_{\alpha})) \quad (\text{factor priors})$$

$$\Lambda(v) = \prod_{\alpha} Z_{\alpha}(v_{\alpha}) \quad (\text{intensity})$$

$$S(v) \sim \mathcal{PO}(S(v); \Lambda(v)) \quad (\text{KL-cost})$$

$$X(v_0) = \sum_{\bar{v}_0} S(v) \quad (\text{observation})$$

$$\hat{X}(v_0) = \sum_{\bar{v}_0} \Lambda(v) \quad (\text{parameter})$$

where the symbols refer to Poisson and Gamma distributions respectively, where:

$$\mathcal{PO}(s; \lambda) = e^{-\lambda} \frac{\lambda^s}{s!} \quad (5)$$

$$\mathcal{G}(z; a, b) = e^{-bz} \frac{z^{a-1} b^a}{\Gamma(a)}. \quad (6)$$

The Gamma prior on the factors are chosen in order to preserve conjugacy. The graphical model for the PLTF framework is depicted in Figure 1. Note that $p(X|S)$ is a degenerate distribution that is defined as follows:

$$p(X|S) = \prod_{v_0} \delta \left(X(v_0) - \sum_{\bar{v}_0} S(v) \right). \quad (7)$$

Here, $\delta(\cdot)$ is the Kronecker delta function where $\delta(x) = 1$ when $x = 0$ and $\delta(x) = 0$ otherwise.

Missing data: To model missing data, we define a 0 – 1 mask array M , the same size as X where $M(v_0) = 1$ ($M(v_0) = 0$) if $X(v_0)$ is observed (missing). Using the mask variables, the missing data is handled smoothly by the following observation model in PLTF:

$$p(X|S)p(S|Z_{1:N}) = \prod_{v_0} \prod_{\bar{v}_0} \{p(X(v_0)|S(v))p(S(v)|Z_{1:N})\}^{M(v_0)} \quad (8)$$

where slight modifications are needed to be done in VB based update equation is shown in section II-B.

B. Fixed Point Update Equation for $PLTF_{KL}$

Here, we recall the generative Probabilistic Latent Tensor Factorization KL model ($PLTF_{KL}$) factor priors with the following fixed point iterative update equation for the component Z_{α} obtained via EM as:

$$Z_{\alpha}(v_{\alpha}) \leftarrow \frac{(A_{\alpha}(v_{\alpha}) - 1) + Z_{\alpha}(v_{\alpha}) \sum_{\bar{v}_0} M(v_0) \frac{\hat{X}(v_0)}{\hat{X}(v_0)} \prod_{\alpha' \neq \alpha} Z_{\alpha'}(v_{\alpha'})}{\frac{A_{\alpha}(v_{\alpha})}{B_{\alpha}(v_{\alpha})} + \sum_{\bar{v}_0} M(v_0) \prod_{\alpha' \neq \alpha} Z_{\alpha'}(v_{\alpha'})} \quad (9)$$

where $\hat{X}(v_0)$ is the model estimate defined as earlier $\hat{X}(v_0) = \sum_{\bar{v}_0} \prod_{\alpha} Z_{\alpha}(v_{\alpha})$. We note that the gamma hyperparameters $A_{\alpha}(v_{\alpha})$ and $B_{\alpha}(v_{\alpha})/A_{\alpha}(v_{\alpha})$ are chosen for computational convenience for sparseness representation such that the distribution has a mean $B_{\alpha}(v_{\alpha})$ and standard deviation $B_{\alpha}(v_{\alpha})/\sqrt{A_{\alpha}(v_{\alpha})}$ and for small $A_{\alpha}(v_{\alpha})$ most of the parameters are forced to be around 0 favoring for a sparse representation [4]. So, equation(9) can be approximated as:

$$Z_{\alpha}(v_{\alpha}) \leftarrow \frac{\sum_{\bar{v}_0} M(v_0) \frac{\hat{X}(v_0)}{\hat{X}(v_0)} \prod_{\alpha' \neq \alpha} Z_{\alpha'}(v_{\alpha'})}{\sum_{\bar{v}_0} M(v_0) \prod_{\alpha' \neq \alpha} Z_{\alpha'}(v_{\alpha'})} \quad (10)$$

Tensor forms via Δ function: We make use of Δ function to make the notation shorter and implementation friendly. A tensor valued $\Delta_{\alpha}^Z(Q)$ function associated with component Z_{α} is defined as follows:

$$\Delta_{\alpha}^Z(Q) = \left[\sum_{\bar{v}_0} \left(Q(v_0) \prod_{\alpha' \neq \alpha} Z_{\alpha'}(v_{\alpha'}) \right) \right] \quad (11)$$

Recall that $\Delta_{\alpha}^Z(Q)$ is an object the same size of Z_{α} while $\Delta_{\alpha}^Z(Q)(v_{\alpha})$ refers to a particular element of $\Delta_{\alpha}^Z(Q)$.

Now, equation(10) can be written into a form that by use of $\Delta_{\alpha}^Z(\cdot)$ as:

$$Z_{\alpha} \leftarrow Z_{\alpha} \circ \Delta_{\alpha}(M \circ X / \hat{X}) / \Delta_{\alpha}(M) \quad (12)$$

where as usual \circ and $/$ stand for element wise multiplication (Hadamard product) and division respectively. We use update equation (12) in the following chapters for PLTF-EM method to compare with the PLTF-VB method.

II. VARIATIONAL BAYES

For a Bayesian point of view, a model is associated with a random variable Θ and interacts with the observed data X simply as $p(\Theta|X) \propto p(X|\Theta)p(\Theta)$. The quantity $p(X|\Theta)$ is called *marginal likelihood* [11] and it is average over the space of the parameters, in our case, S and Z as [4].

$$p(X|\Theta) = \int_Z dZ \sum_S p(X|S, Z, \Theta) p(S, Z|\Theta) \quad (13)$$

On the other hand, computation of this integral is itself a difficult task that requires averaging on several models and parameters. There are several approximation methods such as sampling or deterministic approximations such as Gaussian approximation. One other approximation method is to bound the log marginal likelihood by using *variational inference* [11], [4], [12] where an approximating distribution q is introduced into the log marginal likelihood equation:

$$\log p(X|\Theta) \geq \int_Z dZ \sum_S q(S|Z) \log \frac{p(X, S, Z|\Theta)}{q(S, Z)} \quad (14)$$

where the bound attains its maximum and becomes equal to the log marginal likelihood whenever $q(S, Z)$ is set as $p(S, Z|X, \Theta)$, that is the exact posterior distribution. However, the posterior is usually intractable, and rather, inducing the approximating distribution becomes easier. Here, the approximating distribution q is chosen such that it assumes no coupling between the hidden variables such that it factorizes into independent distributions as $q(S, Z) = q(S)q(Z)$. As exact computation is intractable, we will resort to standard variational Bayes approximations [11], [12]. The interesting result is that we get a belief propagation algorithm for marginal intensity fields rather than marginal probabilities.

A. Variational Update Equations for $PLTF_{KL}$

Here, we formulate the fixed point update equation for the update of the factor Z_α as an expectation of the approximated posterior distribution [13]. Approximation for posterior distribution $q(Z)$ is identified as the gamma distribution with the following parameters:

$$Z_\alpha(v_\alpha) \sim \mathcal{G}(Z_\alpha(v_\alpha); C_\alpha(v_\alpha), D_\alpha(v_\alpha)) \quad (15)$$

where the shape and scale parameters are:

$$C_\alpha(v_\alpha) = A_\alpha(v_\alpha) + \sum_{\bar{v}_\alpha} \frac{X(v_0)}{\hat{X}_L(v_0)} \prod_{\alpha'} L_\alpha(v_\alpha) \quad (16)$$

$$D_\alpha(v_\alpha) = \left(\frac{A_\alpha(v_\alpha)}{B_\alpha(v_\alpha)} + \sum_{\bar{v}_\alpha} \prod_{\alpha' \neq \alpha} \langle Z_{\alpha'}(v_{\alpha'}) \rangle \right)^{-1} \quad (17)$$

Hence the expectation of the factor Z_α is identified as the mean of the gamma distribution and given in the iterative

fixed point update equation obtained via variational Bayes:

$$\begin{aligned} \langle Z_\alpha(v_\alpha) \rangle &= C_\alpha(v_\alpha) D_\alpha(v_\alpha) \\ &= \frac{A_\alpha(v_\alpha) + L_\alpha(v_\alpha) \sum_{\bar{v}_\alpha} \frac{X(v_0)}{\hat{X}_L(v_0)} \prod_{\alpha' \neq \alpha} L_{\alpha'}(v_{\alpha'})}{\frac{A_\alpha(v_\alpha)}{B_\alpha(v_\alpha)} + \sum_{\bar{v}_\alpha} \prod_{\alpha' \neq \alpha} E_{\alpha'}(v_{\alpha'})} \end{aligned} \quad (18)$$

$E_\alpha(v_\alpha)$ and $L_\alpha(v_\alpha)$ (L due to ‘Log’) are two forms of expectations of $Z_\alpha(v_\alpha)$ while $\hat{X}_E(v_0)$ and $\hat{X}_L(v_0)$ are model outputs generated by the components $E_\alpha(v_\alpha)$ and $L_\alpha(v_\alpha)$. While \hat{X}_E is not being used in Equation(19) we define it here, in addition to \hat{X}_L , (and use it later on) since \hat{X}_E has the same shape as \hat{X}_L . Indeed \hat{X}_E and \hat{X}_L can be regarded as different ‘views’ of \hat{X} since they have the same shape (dimensions) as \hat{X} and their computations are done via the same matrix primitives as \hat{X} . Here:

$$E_\alpha(v_\alpha) = \langle Z_\alpha(v_\alpha) \rangle = C_\alpha(v_\alpha) D_\alpha(v_\alpha) \quad (20)$$

$$L_\alpha(v_\alpha) = \exp(\langle \log Z_\alpha(v_\alpha) \rangle) = \exp(\psi(C_\alpha(v_\alpha))) D_\alpha(v_\alpha) \quad (21)$$

$$\hat{X}_E(v_0) = \sum_{\bar{v}_0} \prod_{\alpha} E_\alpha(v_\alpha) \quad (22)$$

$$\hat{X}_L(v_0) = \sum_{\bar{v}_0} \prod_{\alpha} L_\alpha(v_\alpha) \quad (23)$$

Note that the VB version of the update equation(19) closely resembles the EM version given in (9). Indeed when the observed values are large, digamma function becomes $\lim_{x \rightarrow \infty} \psi(x)/\log(x) = 1$, and this, in turn, gives $L_\alpha(v_\alpha) \simeq E_\alpha(v_\alpha)$ and $\hat{X}_L(v_0) \simeq \hat{X}_E(v_0)$.

B. Variational Bound and Sufficient Statistics

The marginal likelihood of the observed data under a tensor factorization model $p(X)$ is often necessary for certain problems such as model selection. We lower bound the marginal likelihood for any arbitrary $PLTF_{KL}$ model based on variational Bayes; while clearly other Bayesian model selection such as MCMC [4] can also be used. To bound the marginal log-likelihood, an approximating distribution $q(S, Z)$ over the hidden structure S and Z is introduced as:

$$\mathcal{L}(\Theta) = \log p(X|\Theta) \geq \int_Z dZ \sum_S q(S, Z) \log \frac{p(X, S, Z|\Theta)}{q(S, Z)} \quad (24)$$

$$= \langle \log p(X, S, Z|\Theta) \rangle_{q(S, Z)} + H[q(S, Z)] = \mathcal{B}_{VB}[q] \quad (25)$$

The bound is tight whenever q equals to the posterior as $q(S, Z) = p(S, Z|X, \Theta)$ but computing the posterior $p(S, Z|X, \Theta)$ is intractable. At this point variational Bayes suggests approximating q . The simplest selection for q from the family of approximating distribution is the one which poses no coupling for the members of the hidden

structure S, Z . That is, we take a factorized approximation $q(S, Z) = q(S)q(Z)$ such that:

$$q(S, Z) = \left(\prod_{v_0} q(S(v_0, *)) \right) \left(\prod_{\alpha} \prod_{v_0} q(Z_{\alpha}(v_{\alpha})) \right) \quad (26)$$

where $*$ symbol in $S(v_0, *)$ is used to indicate the slice of the array. That is $S(v_0, *)$ is the slice of the latent tensor S as the observed variables in configurations v_0 are being fixed. Then, we have:

$$q_{S(v_0, *)}^{(n+1)} \propto \exp \left(\langle \log p(X, S, Z | \Theta) \rangle_{q^{(n)} / q_{S(v_0, *)}} \right) \quad (27)$$

$$q_{Z_{\alpha}(v_{\alpha})}^{(n+1)} \propto \exp \left(\langle \log p(X, S, Z | \Theta) \rangle_{q^{(n+1)} / q_{Z_{\alpha}(v_{\alpha})}} \right) \quad (28)$$

where the superscript (n) indicates the iteration index. This iteration monotonically improves the individual factors of the q distribution, that is, $\mathcal{B}[q^{(n)}] \leq \mathcal{B}[q^{(n+1)}]$ for $n = 1, 2, \dots$ given an initialization $q^{(0)}$.

First, we start with formulating the approximating distribution $q(S)$. When we expand the log and drop $\log P(Z | \Theta)$ and all other irrelevant S terms $q_{S(v_0, *)}$ we end up with:

$$q_{S(v_0, *)} \propto \exp \left(\langle \log p(X | S) + \log p(S | Z) \rangle_{q / q_{S(v_0, *)}} \right) \quad (29)$$

$$\propto \exp \left(\sum_{\bar{v}_0} \left(S(v) \langle \log \prod_{\alpha} Z_{\alpha}(v_{\alpha}) \rangle - \log \Gamma(S(v) + 1) \right) + \log \delta \left(X(v_0) - \sum_{\bar{v}_0} S(v) \right) \right) \quad (30)$$

$$\propto \exp \left(\sum_{\bar{v}_0} \left(S(v) \sum_{\alpha} \log \langle Z_{\alpha}(v_{\alpha}) \rangle - \log \Gamma(S(v) + 1) \right) + \delta \left(X(v_0) - \sum_{\bar{v}_0} S(v) \right) \right) \quad (31)$$

Exactly, the slice $S(v_0, *)$ is sampled from the multinomial distribution as $X(v_0)$ is the total number of observations. Here, s is a vector of a priori independent Poisson random variables s_i . λ is intensity vector conditioned on the sum $x = \sum_i s_i$ and x is multinomial distributed with cell probabilities $p = \lambda / \sum_i \lambda_i$. The joint posterior density of s is denoted by $\mathcal{M}(s; x, p)$. Finally we obtain the approximating distribution as:

$$q_{S(v_0, *)} \sim \mathcal{M}(S(v_0, *), X(v_0), P(v_0, *)) \quad (32)$$

Then, the cell probabilities and sufficient statistics for $q_{S(v_0, *)}$ are:

$$P(v) = \frac{\exp(\sum_{\alpha} \langle \log Z_{\alpha}(v_{\alpha}) \rangle)}{\sum_{\bar{v}_0} \exp(\sum_{\alpha} \langle \log Z_{\alpha}(v_{\alpha}) \rangle)} \quad (33)$$

$$\langle S(v) \rangle = X(v_0)P(v) \quad (34)$$

Now, we turn to formulating $q(Z)$. The distribution $q_{Z_{\alpha}(v_{\alpha})}$ is obtained similarly. After we expand the log and drop irrelevant terms, it becomes proportional to:

$$q_{Z_{\alpha}(v_{\alpha})} \propto \exp \left(\langle \log p(S | Z) + \log p(Z | \Theta) \rangle_{q / q_{Z_{\alpha}(v_{\alpha})}} \right) \quad (35)$$

$$\propto \log Z_{\alpha}(v_{\alpha}) \left(A_{\alpha}(v_{\alpha}) - 1 + \sum_{\bar{v}_{\alpha}} \langle S(v) \rangle \right) - Z_{\alpha}(v_{\alpha}) \left(\frac{A_{\alpha}(v_{\alpha})}{B_{\alpha}(v_{\alpha})} + \sum_{\bar{v}_{\alpha}} \prod_{\alpha' \neq \alpha} \langle Z_{\alpha'}(v_{\alpha'}) \rangle \right) \quad (36)$$

which is the distribution

$$q_{Z_{\alpha}(v_{\alpha})} \sim \mathcal{G}(C_{\alpha}(v_{\alpha}), D_{\alpha}(v_{\alpha})) \quad (37)$$

where the shape and scale parameters for $q_{Z_{\alpha}(v_{\alpha})}$ are given in equation (16) and equation (17).

Finally, sufficient statistics are obtained by the definition of the gamma distribution as follows:

$$E_{\alpha}(v_{\alpha}) = \langle Z_{\alpha}(v_{\alpha}) \rangle = C_{\alpha}(v_{\alpha})D_{\alpha}(v_{\alpha}) \quad (38)$$

$$L_{\alpha}(v_{\alpha}) = \exp(\langle \log Z_{\alpha}(v_{\alpha}) \rangle) = \exp(\psi(C_{\alpha}(v_{\alpha})))D_{\alpha}(v_{\alpha}) \quad (39)$$

Handling Missing Data: Here, slight modifications are needed in the VB-based update equation. We start with the modification on the full joint. Priors are not part of the observation model so they are not affected. The first two terms of $\langle \log p(X, S, Z | \Theta) \rangle_{q(S, Z)}$ become:

$$\sum_{v_0} M(v_0) \left\langle \log \delta \left(X(v_0) - \sum_{\bar{v}_0} S(v) \right) \right\rangle + \sum_{v_0} M(v_0) \left\langle \langle S(v) \rangle \left\langle \log \prod_{\alpha} Z_{\alpha}(v_{\alpha}) \right\rangle - \prod_{\alpha} \langle Z_{\alpha}(v_{\alpha}) \rangle - \langle \log \Gamma(S(v) + 1) \rangle \right\rangle \dots \quad (40)$$

and this results in the following:

$$q_{Z_{\alpha}(v_{\alpha})} \propto \log \langle Z_{\alpha}(v_{\alpha}) \rangle \left(A_{\alpha}(v_{\alpha}) - 1 + \sum_{\bar{v}_{\alpha}} M(v_0) \langle S(v) \rangle \right) - \langle Z_{\alpha}(v_{\alpha}) \rangle \left(\frac{A_{\alpha}(v_{\alpha})}{B_{\alpha}(v_{\alpha})} + \sum_{\bar{v}_{\alpha}} M(v_0) \prod_{\alpha' \neq \alpha} \langle Z_{\alpha'}(v_{\alpha'}) \rangle \right) \quad (41)$$

$$\propto \mathcal{G}(C_{\alpha}(v_{\alpha}), D_{\alpha}(v_{\alpha})) \quad (42)$$

This modifies the gamma parameters for $q(Z)$ given in equations (16) and (17) to include the mask $M(v_0)$ as

follows:

$$C_\alpha(v_\alpha) = A_\alpha(v_\alpha) + \sum_{\bar{v}_\alpha} M(v_0) \langle S(v) \rangle \quad (43)$$

$$D_\alpha(v_\alpha) = \left(\frac{A_\alpha(v_\alpha)}{B_\alpha(v_\alpha)} + \sum_{\bar{v}_\alpha} M(v_0) \prod_{\alpha' \neq \alpha} \langle Z_{\alpha'}(v_{\alpha'}) \rangle \right)^{-1} \quad (44)$$

The other terms are not affected since mask matrix is already in the definition of $C_\alpha(v_\alpha)$ and $D_\alpha(v_\alpha)$. \hat{X}_E and \hat{X}_L are already defined in terms of $C_\alpha(v_\alpha)$ and $D_\alpha(v_\alpha)$. Moreover, $A_\alpha(v_\alpha)$ and $B_\alpha(v_\alpha)$ are priors and not part of the observation model.

Now, for $C_\alpha(v_\alpha)$, we need to find out $\sum_{\bar{v}_\alpha} \langle S(v) \rangle$, which can be written as:

$$\sum_{\bar{v}_\alpha} \langle S(v) \rangle = \sum_{\bar{v}_\alpha} X(v_0) p(v) = \sum_{\bar{v}_\alpha} \frac{X(v_0)}{\hat{X}_L(v_0)} \prod_{\alpha} L_\alpha(v_\alpha) \quad (45)$$

$$= L_\alpha(v_\alpha) \sum_{\bar{v}_\alpha} \frac{X(v_0)}{\hat{X}_L(v_0)} \prod_{\alpha' \neq \alpha} L_{\alpha'}(v_{\alpha'}) \quad (46)$$

After consideration of the missing data for our approach, C_α and D_α can be written using the $\Delta_\alpha^E(\cdot)$ and $\Delta_\alpha^L(\cdot)$ as:

$$C_\alpha = A_\alpha + L_\alpha \circ \Delta_\alpha^L(M \circ X / \hat{X}_L) \quad (47)$$

$$D_\alpha = \left(\frac{A_\alpha}{B_\alpha} + \Delta_\alpha^E(M) \right)^{-1} \quad (48)$$

that, in turn, since $\langle Z_\alpha \rangle$ is $C_\alpha \circ D_\alpha$, E_α and L_α the sufficient statistics for $q(Z_\alpha)$ become:

$$\langle Z_\alpha \rangle = E_\alpha \leftarrow \frac{A_\alpha + L_\alpha \circ \Delta_\alpha^L(M \circ X / \hat{X}_L)}{\frac{A_\alpha}{B_\alpha} + \Delta_\alpha^E(M)} \quad (49)$$

$$\exp(\log(Z_\alpha)) = L_\alpha \leftarrow \exp(\psi(C_\alpha)) \circ D_\alpha \quad (50)$$

After straightforward substitutions, we obtain the variational probabilistic latent tensor factorization algorithm, that can compactly be expressed as in Algorithm 1.

III. EXPERIMENTS AND RESULTS

In this section, we demonstrate the use of the proposed variational Bayesian PLTF (PLTF-VB) for model selection and missing link prediction. First, we study model selection on synthetic datasets and show that the proposed approach can accurately determine the number of components in a CP model. We also show the performance of PLTF-VB for model selection on a real data set, i.e., the UCLAF [14]. Furthermore, on the UCLAF dataset, we study the missing link prediction problem and compare the performance of the proposed variational Bayesian PLTF (PLTF-VB) with the standard PLTF (PLTF-EM) in terms of missing link prediction recovery. For the experiments we use the algorithm that implements variational fixed point update equation given in

Algorithm 1 Variational Inference for PLTF (PLTF-VB)

Input: X (observation), M (mask array), A and B (priors)

Output: E (expected value of factors), and B (bound)

Here $N = |\alpha|$

for $\alpha = 1 \dots N$ **do**

$L_\alpha \sim \mathcal{G}(A_\alpha, B_\alpha / A_\alpha)$

$E_\alpha \sim \mathcal{G}(A_\alpha, B_\alpha / A_\alpha)$

end for

Main loop

for $epoch = 1 \dots MAXITER$ **do**

Compute \hat{X}_L and \hat{X}_E

$\hat{X}_L(v_0) = \sum_{\bar{v}_0} \prod_{\alpha} L_\alpha(v_\alpha)$

Computation for \hat{X}_E is similar and is omitted

for $\alpha = 1 \dots N$ **do**

$C_\alpha = A_\alpha + L_\alpha \circ \Delta_\alpha^L(M \circ X / \hat{X}_L)$

$D_\alpha = 1 / ((A_\alpha / B_\alpha) + \Delta_\alpha^E(M))$

$E_\alpha = C_\alpha \circ D_\alpha$

end for

for $\alpha = 1 \dots N$ **do**

$L_\alpha = \exp(\psi(C_\alpha)) \circ D_\alpha$

end for

end for

panel Algorithm 1 and we use the equation given in (25) for variational bound computation.

Data: As the real data, we use the UCLAF dataset¹ [14] extracted from the GPS data that include information of three types of entities: user, location and activity. The relations between the user-location-activity triplets are used to construct a three-way tensor X . In tensor X , an entry $X(i, j, k)$ indicates the frequency of a user i visiting location j and doing activity k there; otherwise, it is 0. Since we address the link prediction problem in this study, we define the user-location-activity tensor X as:

$$X(i, j, k) = \begin{cases} 1 & \text{if user } i \text{ visits location } j \text{ and} \\ & \text{performs activity } k \text{ there,} \\ 0 & \text{otherwise.} \end{cases} \quad (51)$$

To construct the dataset, the raw GPS points were clustered into 168 meaningful locations and the user comments attached to the GPS data were manually parsed into activity annotations for the 168 locations. Consequently, this dataset consists of 164 users, 168 locations and 5 different types of activities, including ‘Food and Drink’, ‘Shopping’, ‘Movies and Shows’, ‘Sports and Exercise’, and ‘Tourism and Amusement’. In this dataset, 18 users have no location and activity information; it means that the slices corresponding to these users are completely missing. Therefore, we have only used the data from the remaining 146 users. So in our experiments, the number of users is $I = 146$, the

¹<http://www.cse.ust.hk/~vincentz/aaai10.uclaf.data.mat>

number of locations $J = 168$ and the number of activities $K = 5$.

Computational Environment: All experiments were performed using MATLAB 2010b on 2.4GHz Core i5 520M processor and 4GB RAM. Timings were performed using MATLAB's tic and toc functions.

A. Model Selection

Using both synthetic datasets and the UCLAF dataset, we assess the performance of our approach in a model selection context where the goal is to determine the cardinality of the latent index r of the CP model $X^{i,j,k} = \sum_r A^{i,r} B^{j,r} C^{k,r}$. We denote the cardinality of an index i as $|i|$. $|r|$ is set to be from 2 to 10 (ignoring 1) incremented by 1 gradually at each run. In the experiments, the iteration number is set to 2000, the shape parameter A and the scale parameter B of the gamma priors are set to be 0.5 and 10 respectively. As an initialization, a number of random initializations, i.e., 10, are used and the best performing one is picked.

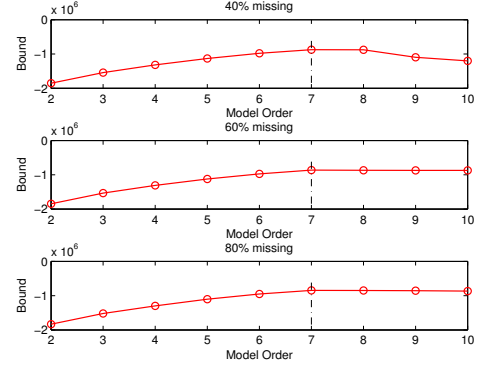
Third-order tensors of different sizes following a CP model are generated. The cardinality of the observed indices $i \times j \times k$ (i.e., the size of the data) are set to $50 \times 50 \times 50$ and $500 \times 500 \times 500$. The cardinality of the latent index r is set to 7 as the true model order. Each run is repeated 10 times and average bound score is plotted in the figures as model order is on the x-axis while bound score given in Equation(24) on the y-axis.

In the first experiment we use the dataset with the size of $50 \times 50 \times 50$ to test model selection process when 40%, 60% and 80% of data is unobserved respectively. What we expect to see on Figure 2(a) is simply that at around true model order of $|r| = 7$ the bound to be the highest to demonstrate that PLTF-VB can find the true model order correctly in all cases of the presence of missing data.

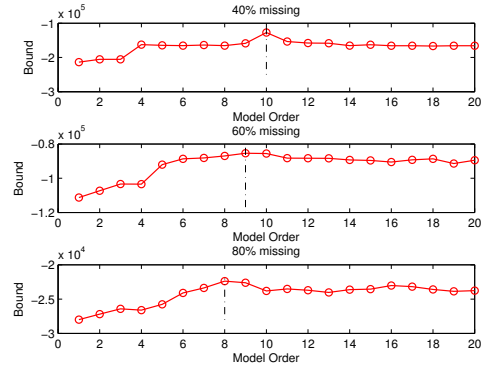
In the second experiment, to illustrate the model order selection performance under missing data case of our model with real data, we use UCLAF dataset. As for the experiment settings, the gamma hyperparameters are set to 0.5 for scale and 10 for shape for all the components and $|r|$ is set to be from 1 to 20. Figure 2(b) shows the performance of our model with 40%, 60% and 80% missing data. We can see from this figure that when the amount of missing data increases the best model order decreases.

B. Scalability

To test the scalability of the PLTF-VB approach, we have done experiments by using the datasets with the size of $50 \times 50 \times 50$ and $500 \times 500 \times 500$. We evaluate the results in terms of both accuracy and speed. Figure 3 shows that PLTF-VB algorithm determines the true model order in both datasets correctly. Then, we run the PLTF-VB algorithm on both datasets ten times and obtain the average solve times. In the $500 \times 500 \times 500$ case, the solve time takes 3974 seconds, approximately 1000 times slower than the



(a) Missing data case with synthetic data



(b) Missing data case with UCLAF

Figure 2. Model order selection using variational bound for CP generated data

$50 \times 50 \times 50$ case (its solve time takes 37 seconds in average), which had 1/1000 times as many variables. In each iteration, the complexity of Algorithm 1 is $O(IJKR)$ where I, J, K are the cardinality of the observed indices and R is the cardinality of the latent index. Consequently, this experiment demonstrates that size of the data and algorithm's complexity are linearly correlated, when one increases the other also increases by the same amount.

C. Hyperparameter Selection

We observe that hyperparameter adaptation is crucial for obtaining good prediction performance. In our simulations, results for PLTF-VB without hyperparameter adaptation were occasionally poorer than the PLTF-EM estimates. We set both shape A and scale B hyperparameters same for all components $Z_{1:3}$. We tried several number of different values for hyperparameters to obtain the best prediction results under missing data case. Figure 4 shows the comparison of three different hyperparameter settings; $A = 0.5, B = 10$, $A = 10, B = 10$ and $A = 100, B = 1$ in terms of link prediction performance. As we can see, we obtain best result

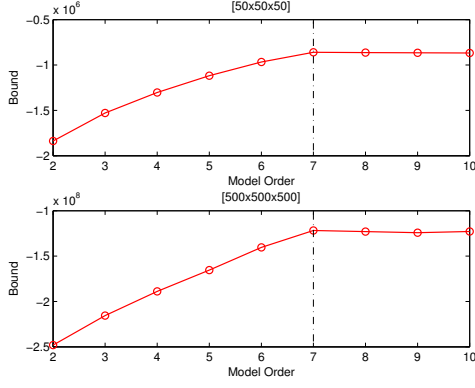


Figure 3. Model order selection on datasets with different sizes

when initialising the shape hyperparameter $A = 0.5$ and scale hyperparameter $B = 10$ for all settings of missing data. So, we use these values of hyperparameter A and B for the following experiments in section III-D. In addition, we obtain that when we set $A < 1$ and $B > 10$, we get better results.

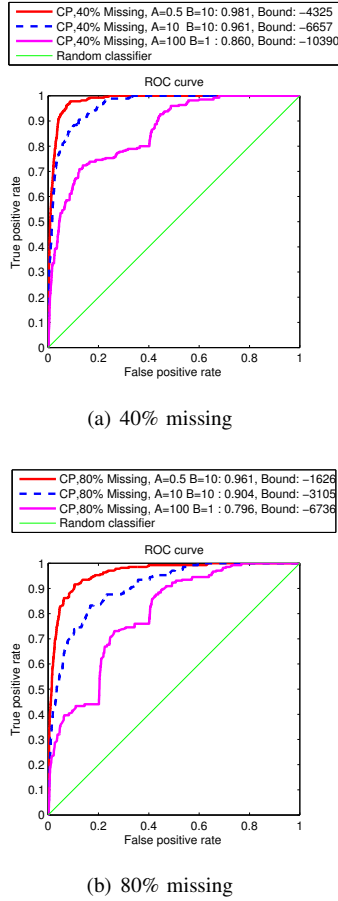


Figure 4. Effect of hyperparameter selection with CP model when $R=2$.

D. Link Prediction

We now compare the standard PLTF, i.e., PLTF-EM, with the proposed variational method, i.e., PLTF-VB, on a missing link prediction task.

1) *Evaluation Metric*: In our experiments, we use Area Under the Receiver Operating Characteristic Curve (AUC) to measure the link prediction performance. Link prediction datasets are characterized by extreme imbalance, i.e., the number of links known to be present is often significantly less than the number of edges known to be absent. This issue motivates the use of AUC as a performance measure since AUC is viewed as a robust measure in the presence of imbalance [15]. The following results show the average link prediction performance of 10 independent runs in terms of AUC.

2) *Results*: We compare the performance of standard and variational approaches of PLTF on both CP and Tucker tensor factorization models at different amounts, i.e., $\{40, 60, 80\}$, of randomly unobserved elements. In these experiments, the incomplete tensor is factorized using either a CP or a Tucker model and the extracted factor matrices are used to construct the full tensor and estimate scores for missing links. For all cases, variational approach outperforms the standard approach clearly. Figure 5 and Figure 6 show the comparison of PLTF-VB and PLTF-EM methods for the CP model given in Equation(3) and the Tucker model given in Equation(4), respectively, when $\{40, 60, 80\}$ of the data is missing. As we can see, the variational methods due to implicit self-regularization effect [16], perform better than the standard methods; in particular, when the percentage of missing data is high. Furthermore, note that the Tucker model outperforms the CP model; because Tucker model is more flexible due to the full core tensor which is helpful for us to explore the structural information embedded in the data.

Moreover, we study the performance of PLTF-EM and PLTF-VB in terms of robustness to model order selection. As model order increases, the prediction performance of PLTF-EM drops. This is as expected since PLTF-EM is prone to overfitting and the increase in model order causes an increase in the number of free parameters that, in turn, enlarges penalty term in PLTF-EM. On the other hand, the prediction performance of the variational approach is not very sensitive to the model order and is immune to overfitting since Bayesian approach alleviates over-fitting by integrating out all model parameters [4]. We compare the prediction performances of PLTF-EM and PLTF-VB methods for the CP tensor model when the component number R is equal to 2 and 20 and for different amounts of missing data, i.e., $\{40, 60, 80\}$ of the data is missing. Figure 7 and Figure 8 demonstrate that when the model order increases, the prediction performance of PLTF-VB approach stays almost same; however, the prediction performance of

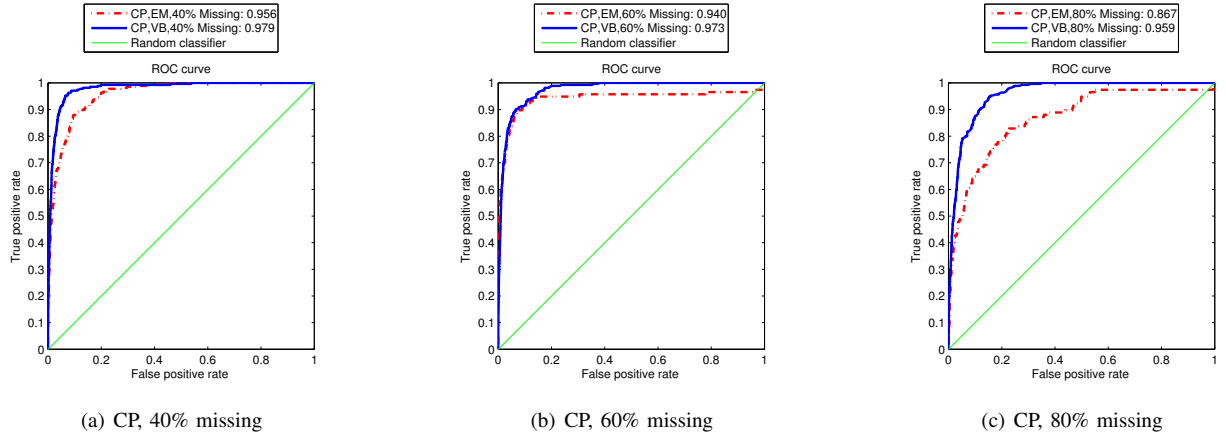


Figure 5. Comparison of PLTF-EM and PLTF-VB methods under missing data case with CP model

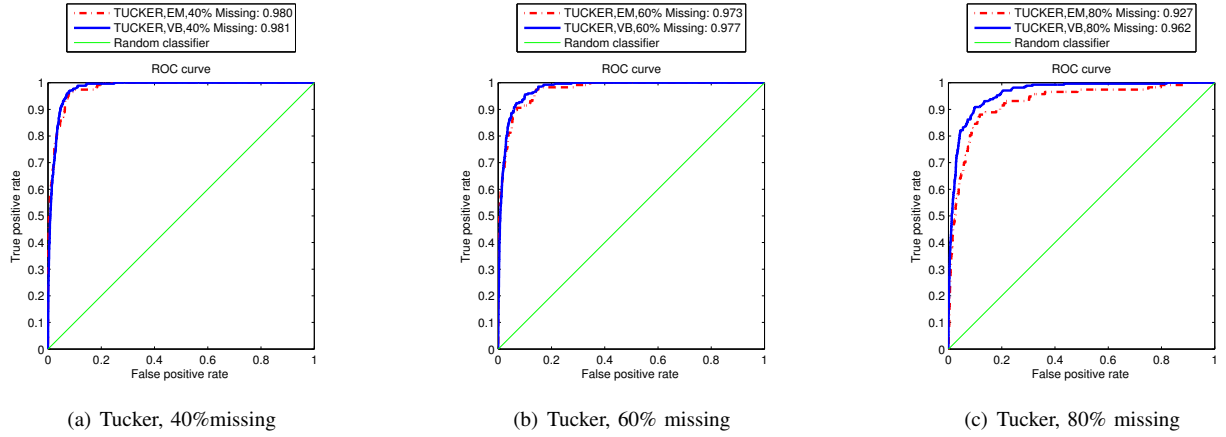


Figure 6. Comparison of PLTF-EM and PLTF-VB methods under missing data case with Tucker model

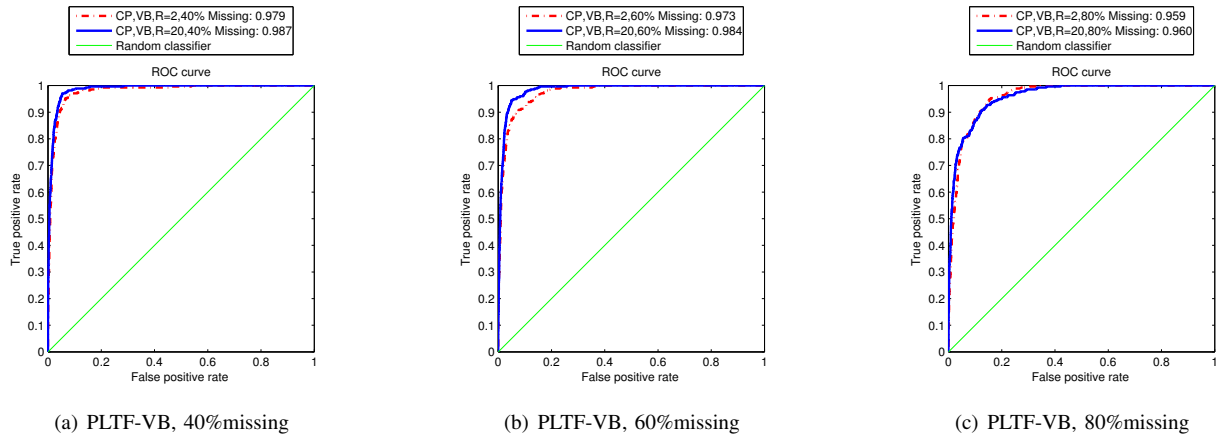


Figure 7. Effect of model order on the performance of PLTF-VB approach for CP model for different amounts of missing data.

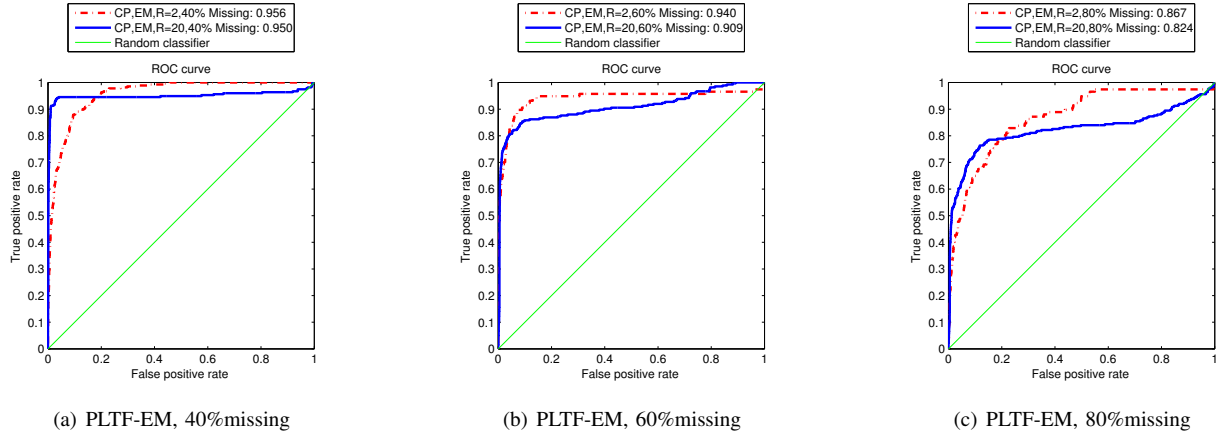


Figure 8. Effect of model order on the performance of PLTF-EM approach for CP model for different amounts of missing data.

PLTF-EM approach declines as expected.

IV. RELATED WORK

In this section, we briefly introduce some of the related work in two categories: Bayesian inference for matrix and tensor factorizations and link prediction.

In order to deal with the variational Bayesian matrix and tensor factorization problem, Ghahramani and Beal [12] provides a method that focus on deriving variational Bayesian learning in a very general form, relating it to EM, motivating parameter-hidden variable factorizations, and the use of conjugate priors. Shan et al. [17] propose probabilistic tensor factorization algorithms, which are naturally applicable to incomplete tensors. First one is parametric probabilistic tensor factorization (PPTF), as well as a variational approximation based algorithm to learn the model and the second one is Bayesian probabilistic tensor factorization (BPTF) which maintains a distribution over all possible parameters by putting a prior on top, instead of picking one best set of model parameters. Cemgil [4] describes a non-negative matrix factorization (NMF) in a statistical framework, with a hierarchical generative model consisting of an observation and a prior component. Starting from this view, he develops full Bayesian inference via variational Bayes or Monte Carlo.

Nakajima et al. [18] propose a global optimal solution to variational Bayesian matrix factorization (VBMF) that can be computed analytically by solving a quartic equation and it is highly advantageous over a popular VBMF algorithm based on iterated conditional modes (ICM), since it can only find a local optimal solution after iterations. Yoo and Choi [19] present a hierarchical Bayesian model for matrix co-factorization in which they derive a variational inference algorithm to approximately compute posterior distributions over factor matrices.

For Bayesian model selection, Sato [20] derives an online version of the variational Bayes algorithm and proves its

convergence by showing that it is a stochastic approximation for finding the maximum of the free energy. By combining sequential model selection procedures, the online variational Bayes algorithm provides a fully online learning method with a model selection mechanism.

We next turn to link prediction studies. Most often, an incomplete set of links is observed and the goal is to predict unobserved links (also referred to as the *missing link prediction* problem), or there is a temporal aspect: snapshots of the set of links up to time t are given and the goal is to predict the links at time $t + 1$ (*temporal link prediction* problem). Matrix and tensor factorization-based methods have recently been studied for temporal link prediction [21]; however, in this paper, we have considered the use of tensor factorizations for the missing link prediction problem. Applications of missing link prediction include predicting links in social networks [22]; predicting the participation of users in events such as email communications and co-authorship [23] and predicting the preferences of users in online retailing [3]. Matrix factorization and tensor factorization-based approaches have proved useful in terms of missing link prediction because missing link prediction is closely related to matrix and tensor completion studies, which have shown that by using a low-rank structure of a data set, it is possible to recover missing entries accurately for matrices [24] and higher-order tensors [25], [26].

V. CONCLUSIONS

In this paper, we have investigated variational inference for PLTF framework with KL cost from a full Bayesian perspective that also handles the missing data naturally. In addition, we develop a practical way without incurring much additional computational cost to PLTF-EM approach for computing the approximation distribution and full conditionals; then, we estimate the model order in terms of marginal likelihood. By maximizing the bound on marginal likelihood, we have a method where all the hyperparameters

can be estimated from data. Our experiments suggest that the variational bound seems to be reasonable approximation to the marginal likelihood and can guide model selection for PLTF.

As a future direction and next step of this work, we aim to extend our variational method in order to be able to make inference on tensor factorization models where multiple observed tensors (X_1, \dots, X_K) can share a set of factors [27]. Factorization of multiple observed tensors simultaneously, alleviates the overfitting better than the standard variational Bayesian matrix factorization and leads to the improved performance [19].

REFERENCES

- [1] Lee, D.D., Seung, H.S.: Learning the parts of objects by non-negative matrix factorization. *Nature* **401** (1999) 788–791
- [2] Paatero, P., Tapper, U.: Positive matrix factorization: A non-negative factor model with optimal utilization of error estimates of data values. *Environmetrics* **5**(2) (1994) 111–126
- [3] Koren, Y., Bell, R., Volinsky, C.: Matrix factorization techniques for recommender systems. *Computer* **42**(8) (2009) 30–37
- [4] Cemgil, A.T.: Bayesian inference in non-negative matrix factorisation models. *Computational Intelligence and Neuroscience* (Article ID 785152) (2009)
- [5] Cichoki, A., Zdunek, R., Phan, A., Amari, S.: *Nonnegative Matrix and Tensor Factorization*. Wiley (2009)
- [6] Yilmaz, Y.K., Cemgil, A.T.: Probabilistic latent tensor factorization. In: *LVA/ICA*. (2010) 346–353
- [7] Tucker, L.R.: Some mathematical notes on three-mode factor analysis. *Psychometrika* **31** (1966) 279–311
- [8] Harshman, R.A.: Foundations of the PARAFAC procedure: Models and conditions for an “explanatory” multi-modal factor analysis. *UCLA working papers in phonetics* **16** (1970) 1–84
- [9] Carroll, J.D., Chang, J.J.: Analysis of individual differences in multidimensional scaling via an N-way generalization of “Eckart-Young” decomposition. *Psychometrika* **35** (1970) 283–319
- [10] Hitchcock, F.L.: The expression of a tensor or a polyadic as a sum of products. *Journal of Mathematics and Physics* **6**(1) (1927) 164–189
- [11] Bishop, C.M.: *Pattern Recognition and Machine Learning* (Information Science and Statistics. Springer (2007)
- [12] Ghahramani, Z., Beal, M.J.: Propagation algorithms for variational bayesian learning. In: *NIPS*. (2000) 507–513
- [13] Yilmaz, Y.K.: *Generalized Tensor Factorization*. PhD thesis, Bogazici University, Istanbul, Turkey (2012)
- [14] Zheng, V.W., Cao, B., Zheng, Y., Xie, X., Yang, Q.: Collaborative filtering meets mobile recommendation: A user-centered approach. In: *AAAI’10: Proceedings of the Twenty-Fourth Conference on Artificial Intelligence*. (2010)
- [15] Stäger, M., Lukowicz, P., Tröster, G.: Dealing with class skew in context recognition. In: *ICDCS Workshops*. (2006) 58
- [16] Nakajima, S., Sugiyama, M.: Implicit regularization in variational bayesian matrix factorization. In: *ICML*. (2010) 815–822
- [17] Shan, H., Banerjee, A., Natarajan, R.: Probabilistic tensor factorization for tensor completion. Technical report, Department of Computer Science and Engineering University of Minnesota (2011)
- [18] Nakajima, S., Sugiyama, M., Tomioka, R.: Global analytic solution for variational bayesian matrix factorization. In: *NIPS*. (2010) 1768–1776
- [19] Yoo, J., Choi, S.: Bayesian matrix co-factorization: Variational algorithm and cramér-rao bound. In: *ECML/PKDD* (3). (2011) 537–552
- [20] aki Sato, M.: Online model selection based on the variational bayes. *Neural Computation* **13**(7) (2001) 1649–1681
- [21] Dunlavy, D.M., Kolda, T.G., Acar, E.: Temporal link prediction using matrix and tensor factorizations. *ACM TKDD* **5**(2) (2011) Article 10
- [22] Clauset, A., Moore, C., Newman, M.E.J.: Hierarchical structure and the prediction of missing links in networks. *Nature* **453** (2008) 98–101
- [23] Getoor, L., Diehl, C.P.: Link mining: a survey. *ACM SIGKDD Explorations Newsletter* **7**(2) (2005) 3–12
- [24] Candès, E.J., Plan, Y.: Matrix completion with noise. *Proceedings of the IEEE* **98** (2010) 925–936
- [25] Acar, E., Dunlavy, D., Kolda, T., Morup, M.: Scalable tensor factorizations for incomplete data. *Chemometrics and Intelligent Laboratory Systems* **106** (2011) 41–56
- [26] Gandy, S., Recht, B., Yamada, I.: Tensor completion and low-n-rank tensor recovery via convex optimization. *Inverse Problems* **27** (2011) 025010
- [27] Yilmaz, Y.K., Cemgil, A.T., Simsekli, U.: Generalised coupled tensor factorisation. In: *NIPS*. (2011)